

Python Machine Learning Labs

Good morning!

We start at 10:00

Data sources

- Excel/CSV/Tabular data sheets
- Databases
- Web

The Setup (1)

You
○
/ \

python
script

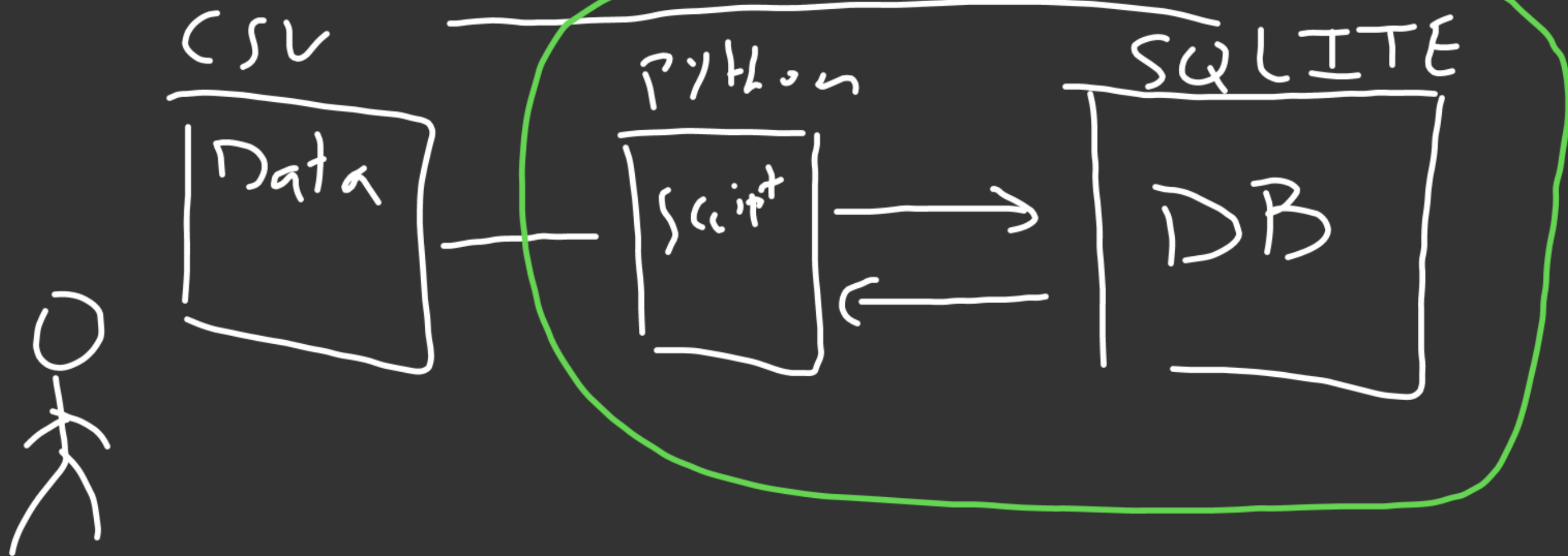
python
code

SQL code



DB
(data)

The Setup (2)



SQLite

- DB is stored on a file
- Lightweight DB
- ✓ - No installation or configuration
- ✗ - Security issues



DB

connection

MySQL Server

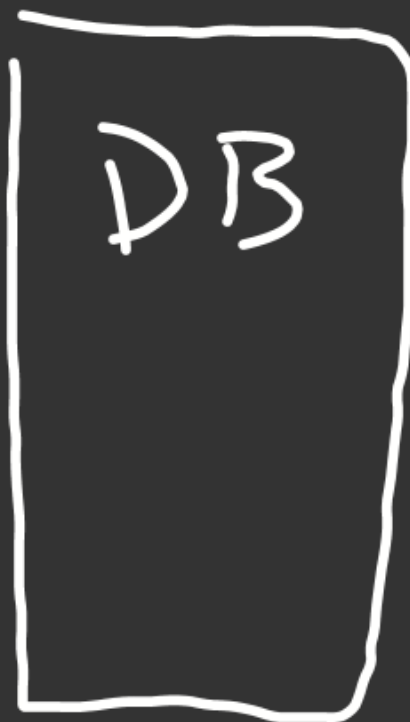
Connect(
dburl
(ip+port)
password

192.17.32.55:55



SQLite

Connect(file path)



Break

Back at 11:15

Feature Engineering

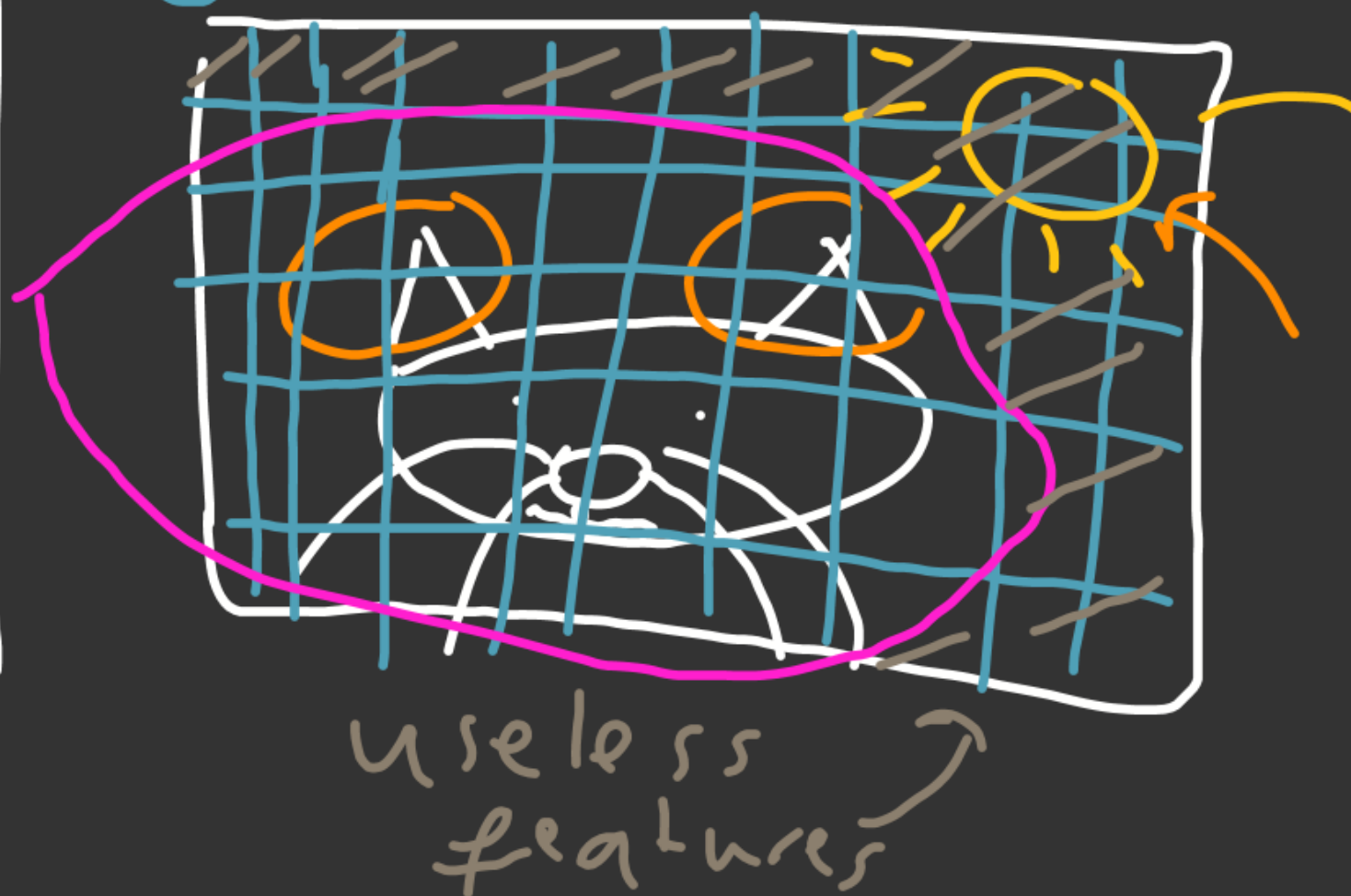
- Fill missing values
 - Drop columns
 - Apply transformations
 - Dimensionality Reduction
- Main Steps

Dimensionality Reduction (1)

Titanic



Feature Cat image



Dimensionality Reduction (2)

Principal Component Analysis (PCA)

n features



2-D features



2-D space

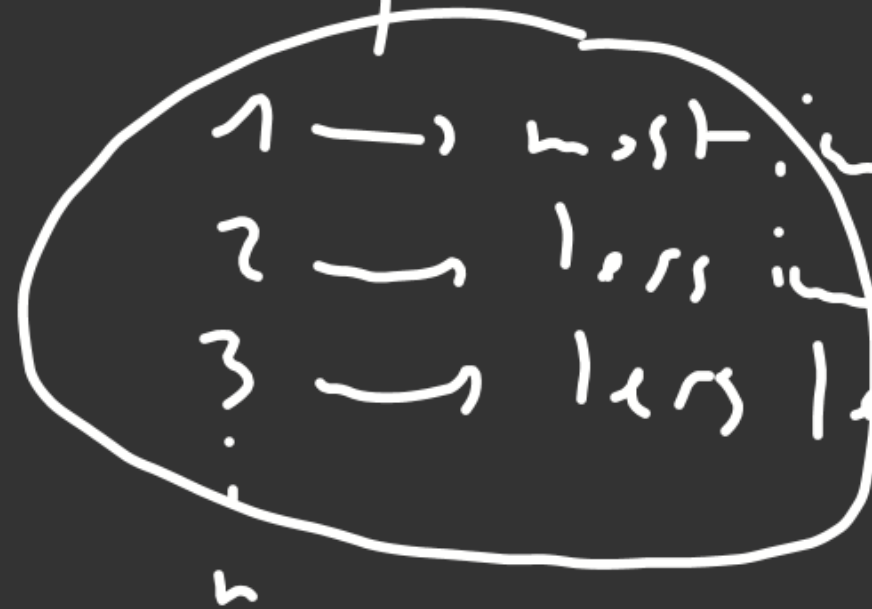
RFE (1)

Feature Selection algorithm
(Recursive Feature Extraction)

n Features



n features



1 → most important

2 → less important

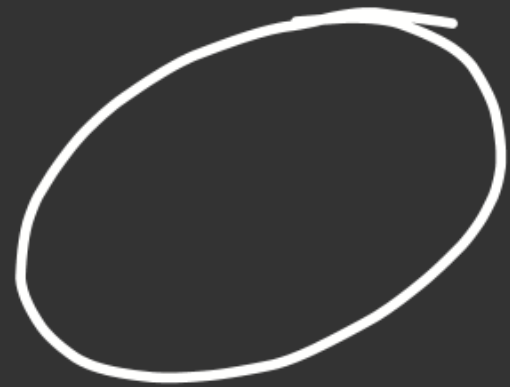
3 → less less important

n

RFE (2)

n features model

+



n-1

features

+ model

RFE

Clustering (1)

Collaborative
Filtering

USER 1

+ shirt?



book
guitar

PROFILE 1

USER 3



— basket ball

PROFILE 2

USER 2



book?

mg

rd playo.

USER 4



— t-shirt

Shunny
Dept

Clustering (2)

Netflix

USER 1

- Pirates of
The Caribbean

- Horror
movies

Secret
Window

USER 2

- Rom - rom

- Hates Jim
Carry

USER 3

- Likes Hugh
Grant

Love Actually
(Hugh Grant +
Rom - Com)

Lunch

Back at 14:00

Filling missing Age (1)

id	Class	Sex	Age
0	1	f	38
1	2	f	25
2	1	f	30
3	1	f	

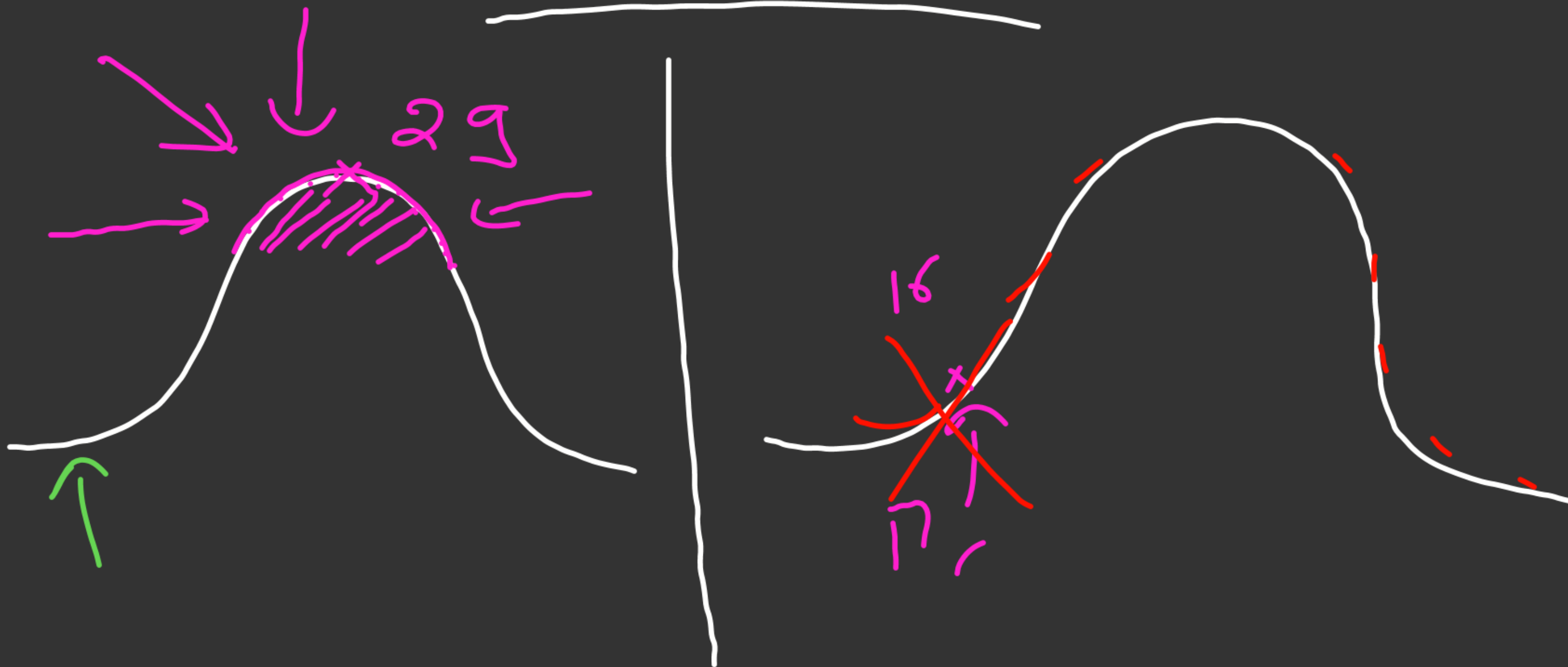
id	Class	Sex	Age
0	1	f	38
2	1	f	30
3	1	f	

Filling Missing Ages (2)

id	relat	sex	Age
0	1	f	38
2	1	f	30
3	1	f	X

[3, ...]

Filling Data



Break

Back at 15:05

Features

Numerical

- Discrete

e.g. 1, 2, 3 ...

- Continuous

e.g. 1.1, 2.1, 1.001, ...

Categorical

- Nominal

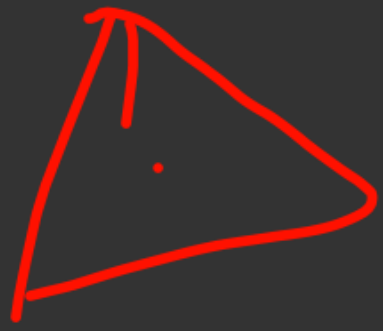
e.g., Japan, USA, ...

- Ordinal

e.g., high, medium, low ...

- Free text

e.g., summary of a book



Encoding (1)

Label

Country

Japan

USA

India

①

②

③

India >> Japan

Pros & Cons

- Numerical ✓
- Good for many values ✓
- Can infer wrong relationships (ordinal) ✗

Leave-one-out
= Represent k values,
with $k-1$ columns

Encoding (2)

Vector (one-hot)

Country

Japan

USA

India

J	U	I
1	0	0
0	1	0
0	0	1

Pros & Cons

- No inference of wrong relationships ✓
- Not good for many value representation ✗
- Memory overload ✗

Encoding (3)

Sex
male
female



m	f
1	0
0	1